

A Statistical Treatment of Bioassay Pour Fractions

Jack B. Barengoltz, Ph.D. and David Hughes (GSFC)

The binomial probability distribution is used to treat the statistics of a microbiological sample that is split into two parts, with only one part evaluated for spore count. One wishes to estimate the total number of spores in the sample based on the counts obtained from the part that is evaluated (pour fraction). Formally, the binomial distribution is recharacterized as a function of the observed counts (successes), with the total number (trials) an unknown. The pour fraction is the probability of success per spore (trial). This distribution must be renormalized in terms of the total number. Finally, the new renormalized distribution is integrated and mathematically inverted to yield the maximum estimate of the total number as a function of a desired level of confidence ($P(<n)=LOC$).

Selected results of the indicated numerical calculations are presented. For $LOC=0.5$, or the likely value, the estimates differ little from the usual calculation: the number of spores counted divided by the pour fraction.

The extension to recovery efficiency corrections is also presented. Now the product of recovery efficiency and pour fraction may be small enough that the likely value may be much larger than the usual calculation: the number of spores divided by that product.

The use of this analysis would not be limited to microbiological data.

BIOASSAY POUR FRACTION STATISTICS

In a surface microbial assay, a spacecraft surface is sampled manually with a damp sterile assay tool (swab or wipe). The sampling tool is placed in a vessel filled with sterile water. The sample is processed by heat shock to kill the microbes that are not spores (inexact but the requirements are based on the procedure). The sample is further treated by sonication to remove the spores from the sampling tool, from the vessel wall, from particles and from each other. Vortexing helps to suspend the single spores in a uniform suspension. Then part of the rinse suspension is poured or streaked onto plates with growth medium. After further microbiological processing, the colonies that form on the plates are counted. This raw data is then used to estimate the spores on the area sampled.

Consider the rinse suspension alone and the pour fraction. One needs to estimate the number of spores in the rinse solution from the count on the plate. The usual calculation is simply to divide the raw counts by the

pour fraction. A more detailed method is possible by considering the pour fraction a random selection of a part of the rinse solution (and a random selection of the total number of spores in the rinse solution). The binomial distribution provides the necessary probability distribution.

$$B(n_0, n, f) = C(n, n_0) f^{n_0} (1-f)^{(n-n_0)}$$

where

$$C(n, n_0) = \frac{n!}{n_0! (n - n_0)!}$$

n_0 is the number of spores counted (successes), n is the true number of spores in the rinse and f is the pour fraction. Note that $n \geq n_0$, of course.

However, in the problem at hand, one requires the probability distribution of n (which is unknown) given that n_0 spores are counted in the pour fraction f . The distribution above represents that of n_0 given known n . This distribution must be normalized so that the sum over n , starting at n_0 , to ∞ is unity.

For all cases except $n_0=0$, the required sum was performed in EXCEL, with the use of the BinomDist function. It may be shown explicitly that the sum for the case of $n_0=0$ is $1/f$. By direct calculation, this is also true for all other values of n_0 . The combination factor $C(n, n_0)$ eliminates any easy proof. The desired distribution is therefore

$$B'(n_0, n, f) = f B(n_0, n, f)$$

Only the distribution form of BinomDist may be used (last argument = FALSE). The integral form of the new distribution is

$$P(< n, n_0, f) = \sum_{n=n_0}^{\infty} B'(n_0, n, f)$$

which must be calculated.

The integral form of BinomDist (last argument =TRUE) is of no use (sum is on n_0 , among other issues).

Figures 1a and 1b show the results of calculations for the differential (probability density) and integral probability distributions $p(n, n_0, f)$ and $P(<n, n_0, f)$, respectively, for $f=0.8$, a typical value for a swab sample. Results for a few values of n_0 are shown. Figures 2a and 2b show the results for $f=0.25$, a typical value for a wipe sample.

Of course for increasing n_0 , the distributions move toward increasing n . The differential distributions also broaden, and the slope of the integral distributions decrease.

This analysis was originally performed to determine the statistics of the raw counts. To that end, the counts n corrected for the pour fraction for a selected level of confidence (LOC) may be calculated with the integral probability $P(<n)$. The results calculated for LOC corresponding to Gaussian mean + 1σ , + 2σ , and + 3σ are shown in Tables 1 and 2. Note that all of this analysis is for integer counts; exact correspondence is impossible. At the same time, the most likely integer value for n , as close to corresponding to $P(<n)=0.5$ as possible is shown for comparison with the current method (not an integer), n_0/f .

Table 1. Statistics for n , $f = 0.8$

measured counts	0	1	3	10
measured counts/ f	0	1.25	3.75	12.5
likely max counts LOC =0.5	0	0-1	3-4	12
max counts at LOC=.84	0-1	1-2	4-5	14
max counts at LOC=.977	1-2	3	6	16-17
max counts at LOC=.9987	3	5	8	20

Table 2. Statistics for n , $f = 0.25$

measured counts	0	1	3	10
measured counts/ f	0	4	12	40
likely max counts LOC =0.5	1-2	5-6	13-14	41-42
max counts at LOC=.84	5-6	11	21	54
max counts at LOC=.977	12	19	31-32	68-69
max counts at LOC=.9987	22	30-31	45	86-87

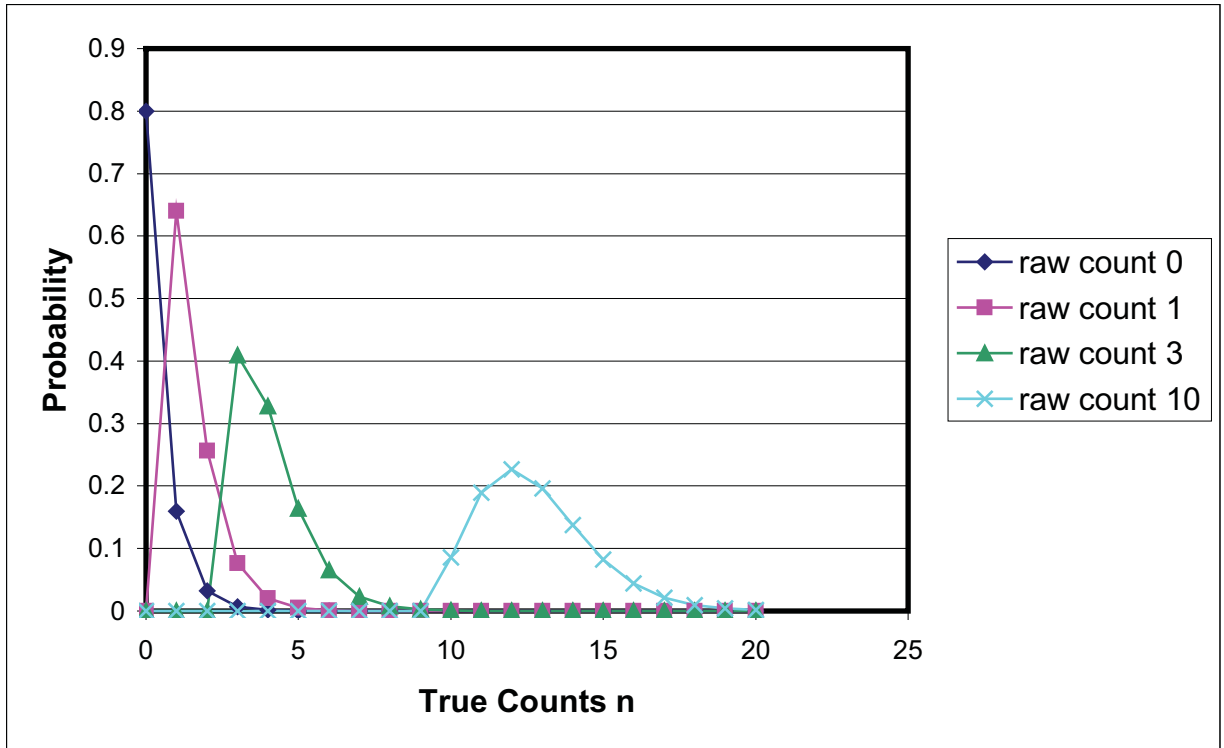


Figure 1a. Probability Density of True Counts n for $f = 0.8$ for indicated raw counts

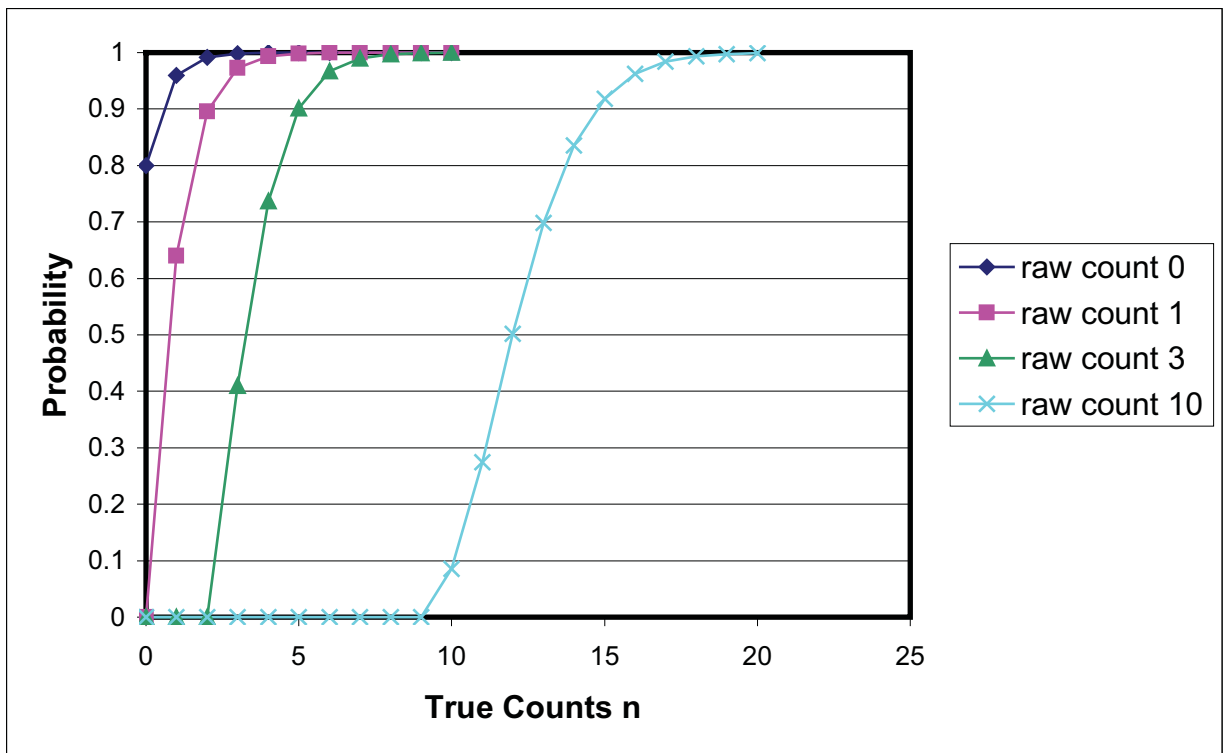


Figure 1b. Probability of True Counts $\leq n$ for $f = 0.8$ for indicated raw counts

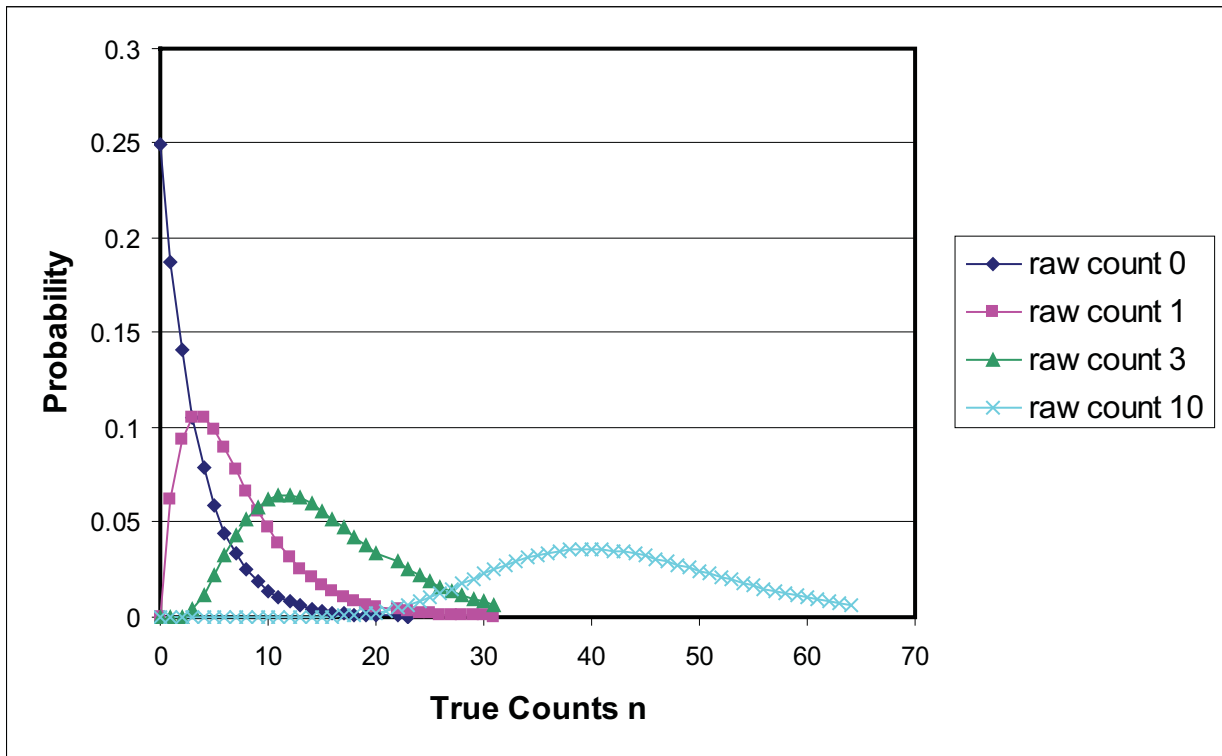


Figure 2a. Probability Density of True Counts n for $f = 0.25$ for indicated raw counts

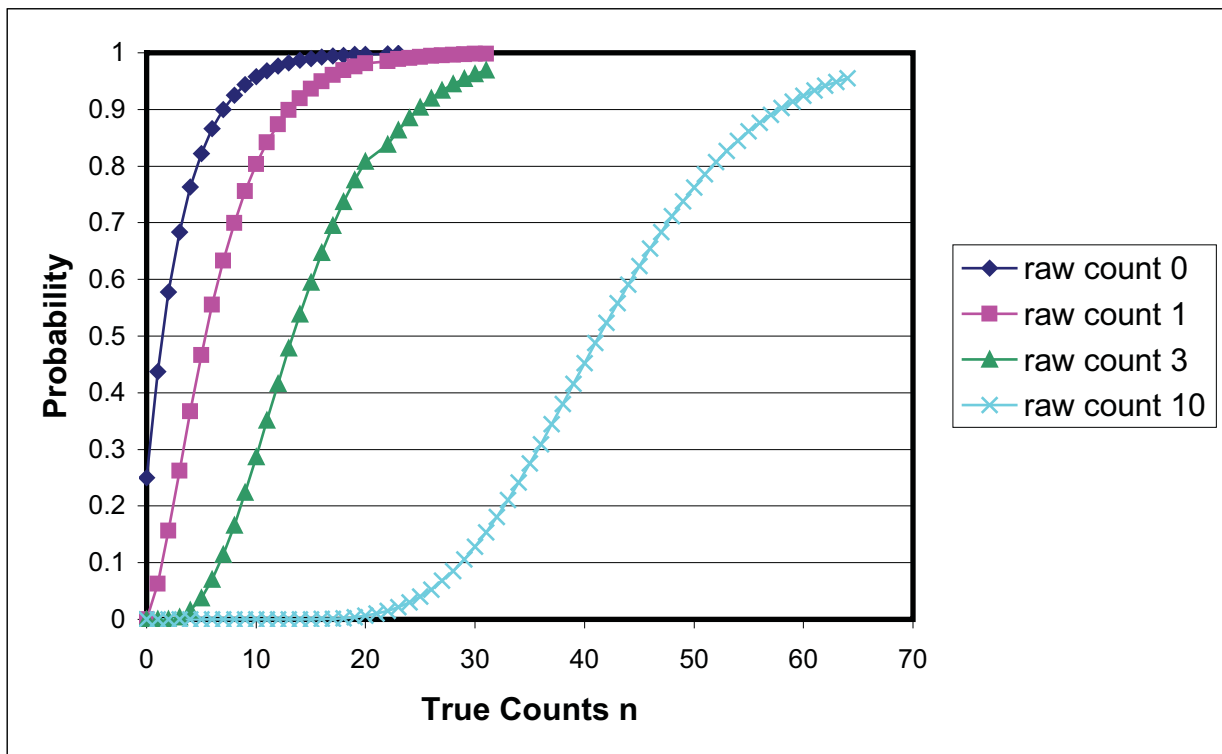


Figure 2b. Probability of True Counts $\leq n$ for $f = 0.25$ for indicated raw counts

It can be seen that for small pour fractions, the increases for higher LOC can be very large.

In addition the (raw) measured counts divided by the pour fraction (n_0/f) is a reasonable correction at $f=0.8$. These values compare favorably with the max counts for $LOC=0.5$. The comparison is good even for the case of zero measured counts, when the correction has no effect, because the probability of max counts=0 is 0.8 (>0.5).

However, for $f=0.25$, n_0/f is a non-conservative estimate of the likely maximum counts. The difference may be as large as 2 counts, which is extremely important for 0 or one raw count cases.

EXTENSION TO RECOVERY EFFICIENCY (e)

This analysis is not limited to the probability that a spore is in the fraction of the rinse solution poured or streaked. If the recovery efficiency that relates the plate counts to the true count on the surface sampled is known, then the same binomial probability applies. Successes are the number of plate counts. The number of spores on surface sampled (not just in rinse) becomes the n true count for which statistics are needed. The recovery efficiency is the product of the probabilities of spores on surface removed, released from the sampling tool, and poured, etc. The recovery efficiency is measured separately with known samples. Since the value is explicitly corrected for pour fraction f , the overall probability of a count from one spore on the surface is ef .

Table 3. Pour Fractions and Recovery Efficiencies

	Pour Fraction, f	Recovery Efficiency, e	Product, ef
ESA swab	0.8	0.45	0.36
ESA wipe	0.25	0.20	0.05
NASA swab	0.8	0.30	0.24
NASA wipe	0.25	tbd	tbd

The results of this analysis for $ef=0.36$ and 0.05 are shown in Tables 4 and 5. The analysis is not changed, but since $ef < f$, the statistics lead to higher estimates for n at same LOC. In particular, for the very small ef for ESA wipes, 0.05, n_0/ef greatly underestimates even the likely max

counts. The results for the ESA swabs are comparable. However, for one or two counts, n_0/ef is too small.

The statistics for the NASA swab are not shown because they are very similar to Table 2, for $f=0.25$.

Table 4. Statistics for n , $ef=0.36$

measured counts	0	1	3	10
measured counts/ f	0	2.8	8.3	27.8
likely max counts LOC =0.5	0-1	3-4	9	28-29
max counts at LOC=.84	3	7	14	36-37
max counts at LOC=.977	7-8	12-13	20-21	46
max counts at LOC=.9987	14	19-20	29	57-58

Table 5. Statistics for n , $ef=0.05$

measured counts	0	1	3	10
measured counts/ f	0	20	60	200
likely max counts LOC =0.5	12-13	32-33	81	221
max counts at LOC=.84	35	63-64	115-116	281-282
max counts at LOC=.977	72-73	110	174	366
max counts at LOC=.9987	128-129	174	248-249	466-467

STATISTICS FOR A SAMPLE SET

A sample set is a group of samples that are assumed to be of the same population. A set of samples from the surface of a single spacecraft hardware item is an example. The usual method is to sum the counts of the group and divide by the product of the pour fraction and the recovery efficiency. (If there are swabs and wipes in the set, then the swab total counts and the wipe total counts are corrected separately and then added.)

The problem is that the solution for the total count statistics for a sample set is not unique. More than one set of estimates at various LOCs of the actual counts on each sample's surface can correspond to the same overall LOC. These cases have different values for the total estimated count.

The simplest approach with this new analysis is to apply the statistics for $LOC=0.50$ for the most likely max total number. Apply the statistics for other LOCs as desired to get a count estimate for the surface area sampled by the set.

One should not use the estimated counts at some LOC for each sample's count and sum over the sample set. This estimate is much less likely to be exceeded than $1-LOC$. In fact the effective LOC of this unlikely arrangement is $1-(1-LOC)^N$ where N is the number of samples in the set. It is very conservative and unrealistic, but is always an upper limit.

CONCLUSIONS

The analysis presented is useful whenever the binomial probability distribution applies, and an estimate of the total number of trials is needed given an observation of the number of successes. Such an analysis requires that each trial is independent and equivalent. For planetary protection that means that the spores do not interact physically for a single sample. For a sample set, it means that the microbiological procedure is well controlled (repeatable), and the sampled areas are equivalent for the set (a judgement on the source and history of the hardware of the sample set).

ACKNOWLEDGEMENT

The authors wish to thank Dr. Catherine Conley, the NASA Planetary Protection Officer, for her financial support of this work. One of us (Barengoltz) wishes to thank also Dr. Gerhard Kminek, the ESA Planetary Protection Officer, for his financial support.